

基于多因素稀疏回归预测模型的商家客流量预测 *

郑增威¹, 杜俊杰^{1,2}, 周燕真^{1,2}, 孙霖^{1†}, 霍梅梅¹

(1. 浙江大学城市学院 智能植物工厂浙江省工程实验室, 杭州 310015; 2. 浙江大学 计算机科学与技术学院, 杭州 310012)

摘要: 针对智能商业平台中的大数据预测问题, 提出一种多因素稀疏回归预测模型。以离散余弦变换为基础, 构建包含多个外部因素(节假日、天气、温度)的字典集, 通过 LASSO 方法定量求解稀疏编码模型中各外部因素的影响。实验对 2 000 个商家的客流量进行预测。实验结果表明, 外部因素不同程度地影响客流量, 在预测模型中叠加外部因素后可以有效提高预测的准确性。同时, 通过与其他方法的对比, 多因素稀疏回归预测模型比 RNN、ARIMA 等模型的预测效果更好。

关键词: 智能商业平台; 客流量预测; 稀疏回归; 多因素分析; 字典学习

中图分类号: TP391 **doi:** 10.19734/j.issn.1001-3695.2018.10.0816

Retail consumer traffic forecasting based on multi-factor sparse regression prediction model

Zheng Zengwei¹, Du Junjie^{1,2}, Zhou Yanzhen^{1,2}, Sun Lin^{1†}, Huo Meimei¹

(1. Intelligent Plant Factory of Zhejiang Province Engineering Laboratory, City College, Zhejiang University, Hangzhou 310015, China; 2. College of Computer Science & Technology, Zhejiang University, Hangzhou 310012, China)

Abstract: This paper proposed a multi-factor sparse regression prediction model aiming to solve the problem of big data prediction in business intelligent platform. Constructed a dictionary containing external factors (holidays, weather, and temperature) based on the discrete cosine transform, and quantitatively solved the influence of external factors in the sparse coding model by LASSO. In experiments, the customer traffics of 2 000 stores were predicted. The experimental results show that the impact of external factors on the store customer traffic are different, and the prediction accuracy can be effectively improved with the combination of external factors in the prediction model. In addition, the method was compared with other forecasting methods. The result shows that multi-factor sparse regression prediction model outperforms than other models such as RNN and ARIMA.

Key words: business intelligent platform; customer traffic prediction; sparse regression; multiple factors analysis; dictionary learning

0 引言

随着移动互联网的快速发展, 人们可以方便地使用手机选择附近的商家进行消费, 这使得零售业务的竞争越来越激烈。零售服务是典型的定制服务, 高效的库存管理^[1,2]是满足客户需求的基础。不准确的消费者流量预测可能导致库存过多或不足, 这将直接影响商家业务的盈利能力和竞争地位。而消费者流量的准确预测可以通过提高连锁经营效率和尽量减少浪费来提高零售商的盈利能力。因此, 对于零售商店来说, 根据准确的客流量预测来制定正确的营销策略尤为重要。

零售业务消费者流量预测是对时间序列的短期预测, 它依赖于历史数据并预测未来的消费者流量。自回归移动平均 (autoregressive integrated moving average model, ARIMA)^[3]模型是时间序列预测中最广泛使用的经典方法。ARIMA 预测框架最早是由 Box 和 Jenkins 开发的, 这个框架包括模型选择、参数估计和模型检验的三个迭代过程。Ramos 等人^[4]在消费者零售销售预测案例研究中比较了 ARIMA 模型和指数平滑的预测效果, 实验结果表明 ARIMA 模型要更优于指数

平滑预测方法。文献[5]使用模糊时间序列模型和季节模型对南京某商场的客流量进行预测, 实验结果表明季节模型要优于模糊时间序列。Liu 等人^[6]在快餐店数据预测案例中的提出了时间序列数据挖掘方法, 并对 Box-Jenkins 时间序列预测方法进行了改进。灰色系统预测^[7]是时间序列的另一种方法, 用于解决有限数据和信息不足的不确定问题。灰色预测与 ARIMA 模型不同, 它对于受不确定因素影响较大的复杂环境预测效果较好, 而且所需的样本数据较小。文献[8]提出一种离散灰色预测模型与人工神经网络混合的智能模型, 实验结果表明该算法可有效地用于时尚销售的即时预测。数据挖掘技术提供了一种将大量数据分解成可应用于时间序列分析的信息的方法。İrem 等人^[9]提出了一种数据挖掘方法来预测零售需求, 他们采用二分组聚类算法对具有相似销售行为的仓库进行分组, 并采用贝叶斯网络获得较好的预测结果。Schneider 等人^[10]提出了一种基于属性回归模型的随机预测方法来预测商品的销售量, 实验结果表明用户评论对于商品的销售有着明显的影响, 该方法具有很好的泛化性与扩展性。文献[11]基于用户消费行为数据提出了一种结合 Huff 模型^[12]

收稿日期: 2018-10-29; 修回日期: 2018-12-16 基金项目: 浙江省自然科学基金资助项目 (LY17F020008)

作者简介: 郑增威 (1969-), 男, 浙江温州人, 教授, 硕导, 博士, 主要研究方向为物联网大数据分析、普适计算; 杜俊杰 (1993-), 男, 浙江金华人, 硕士研究生, 主要研究方向为数据挖掘、机器学习; 周燕真 (1996-), 男, 江西吉安人, 硕士研究生, 主要研究方向为机器学习、数据挖掘; 孙霖 (1979-), 男 (通信作者), 浙江杭州人, 副教授, 博士, 主要研究方向为机器学习、普适计算 (sunl@zucc.edu.cn); 霍梅梅 (1977-), 女, 山东德州人, 副教授, 硕士, 主要研究方向为无线传感网络、嵌入式系统应用开发。

和 Monte Carlo 模拟^[13]的预测方法, Huff 模型能很好反映商店与客户之间的关系, 但是对一些市场潜在规则以及用户行为却无法很好表示, 而 Monte Carlo 模拟能克服这个缺点。因此, 该方法的基本思想是挖掘两种模型之间的关系来进一步提高预测精度。

综上所述, 目前时间序列预测基本都是仅仅依靠历史数据来预测未来一段时间的趋势。但是时序数据不仅仅与历史客流量相关, 特别是商家客流量, 还可能与未来一段时间的节假日、天气、温度一系列外部因素有关。而这些因素常常耦合在一起, 因此商家客流量预测成为一个难以建立有效数学模型的复杂的, 高度不确定的非线性波动系统。使用传统方法难以对商店客流量建立有效的预测模型。本文综合考虑了多种影响商家客流量的外部因素, 提出了一种结合节假日、温度、天气情况的多因素稀疏回归预测模型, 将外部因素按稀疏编码系数叠加以增强预测准确性, 该模型在 2000 个商家上进行客流量预测, 同时分析了这些因素对预测精度的影响。实验结果表明, 附加这些外部因素后可以有效提高客流量预测的准确性。

1 多因素稀疏回归模型

1.1 模型框架

本文构建的多因素稀疏回归模型的预测流程如图 1 所示。首先, 对训练数据进行预处理, 去除噪声数据以及数据标准化; 接着, 使用离散余弦变换(discrete cosine transform, DCT)^[14], 克罗内克函数以及外部因素(节假日、温度和天气)构建一个过完备的多因素字典; 然后, 使用构建的字典对训练数据进行稀疏分解, 求解稀疏系数 α ; 最后, 根据稀疏系数和字典进行未来一段时间的客流量的预测。

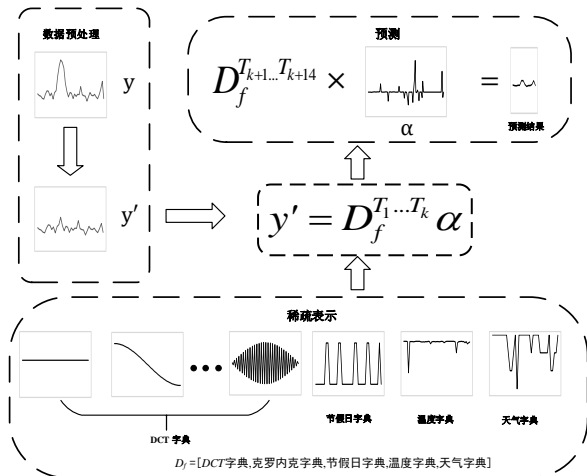


图 1 多因素稀疏回归预测模型框架

Fig. 1 Framework of multi-factors sparse regression prediction model

1.2 客流量数据预处理

因为商店的促销、暂停营业或者其他原因, 一些商店客流时序数据中包含很多噪声, 这些噪声数据或远大于正常数据, 或远小于正常数据。在本文中, 例用分位数特征来过滤每个店铺数据中的异常值。处理过程如图 2 所示。

在本文中使用式(1)对训练数据进行标准化操作。

$$y'' = \frac{y' - \text{mean}(y')}{\text{std}(y')} \quad (1)$$

其中: y' 表示去噪之后的训练数据; $\text{mean}(y')$ 和 $\text{std}(y')$ 分别表示 y' 的均值和标准差; y'' 表示标准化的结果。

值得注意的是, 图 2 表示的是一次预测中数据的预测处理过程, 在对商家多次的预测过程中, 每次处理的都是训练

数据, 而不是商家所有的数据。因此, 数据预处理是一个不断循环的过程。

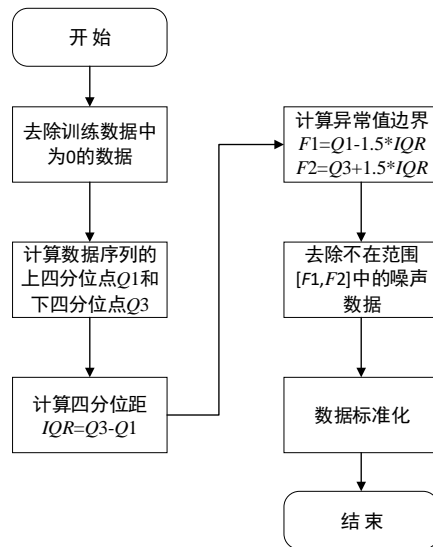


图 2 数据预处理流程

Fig. 2 Process of data preprocessing

图 3 是数据去噪的示意图, 其中图 3(a)是原始数据, 可以看到有明显的异常值; (b)是去除异常值之后的数据。

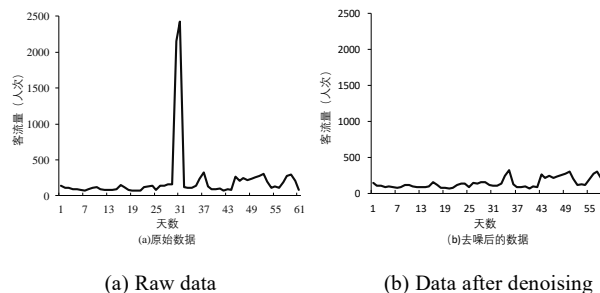


图 3 数据去噪示意图

Fig. 3 Example of data denoising

1.3 多因素字典

字典对后续稀疏系数向量的求解以及客流量的预测十分重要, 目前字典构建主要有基于分析和基于学习^[15]两种方法。在本文中, 由于外界因素对客流量的影响, 所以使用基于分析的方法, 以离散余弦变换为基础, 综合节假日因素、温度因素、天气因素, 构建一个过完备多因素字典 D_f 用于稀疏编码。

1.3.1 过完备 DCT 字典

离散余弦变换 (DCT)^[14]是一种变换压缩方法, 在信号、图像处理中被广泛使用, DCT 使用余弦函数来表示信号量, 它有几种变体。在本文中使用式(2)中所示的正交 DCT-II 来构建大小为 $N \times N$ 的 DCT 字典。

$$\varphi_i(j) = \begin{cases} \frac{1}{\sqrt{N}} & i=0 \\ \sqrt{\frac{2}{N}} \cos\left(\frac{(2j+1)\pi}{2N}i\right) & i=1, 2, \dots, N-1 \end{cases} \quad (2)$$

同时, 为了减少过拟合, 本文引入式(3)所示的 Kronecker Delta 函数^[16]构建另一个大小为 $N \times N$ 克罗内克字典, 并与上述 DCT 字典组成新的过完备字典。

$$K(i, j) = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases} \quad (3)$$

式(2)和(3)中 $i, j \in [0, N-1]$ 表示 DCT 字典和克罗内克字典中第 i 列、第 j 行的元素。根据上述定义, 本文最终组建的

大小为 $N \times 2N$ 的字典 D 如式(4)所示, 前 N 列是 DCT 产生的子字典, 最后 N 列是由 Kronecker Delta 函数生成的子字典。

$$D = \begin{bmatrix} \frac{1}{\sqrt{N}} & \sqrt{\frac{2}{N}} \cos \frac{\pi}{2N} & \dots & \sqrt{\frac{2}{N}} \cos \frac{(N-1)\pi}{2N} & 1 & 0 & \dots & 0 \\ \frac{1}{\sqrt{N}} & \sqrt{\frac{2}{N}} \cos \frac{3\pi}{2N} & \dots & \sqrt{\frac{2}{N}} \cos \frac{3(N-1)\pi}{2N} & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\sqrt{N}} & \sqrt{\frac{2}{N}} \cos \frac{2(N-1)\pi}{2N} & \dots & \sqrt{\frac{2}{N}} \cos \frac{2(N-1)(N-1)\pi}{2N} & 0 & 0 & \dots & 1 \end{bmatrix} \quad (4)$$

1.3.2 节假日字典

节假日因素与人们的出行有很大的关系, 间接地会对商店的客流量产生影响。与工作日相比, 人们更有可能在节假日外出游玩。因此, 一些商店的消费者流量在节假日将会明显提高。另一方面, 以快餐店为典型的这类商店客流量可能反而会在节假日下降。无论何种情况, 节假日与客流量之间都存在关联关系。在本文中构建了节假日字典, 使用 **1** 表示节假日, **0** 为工作日。另外, 为了更好地契合实际情况, 本文根据中国法定假日重新调整节假日字典的值。图 4 表示一个节假日字典例子。

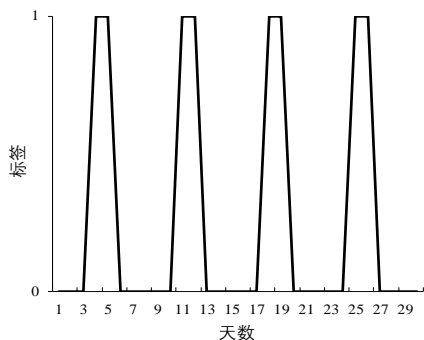


图 4 节假日字典例子

Fig. 4 Example of holiday dictionary

1.3.3 温度字典

对于温度而言, 当温度适宜时, 人们更愿意出行。但是几度的温差并不会对人的活动造成很大的影响。考虑到中国相对稳定的天气条件, 极端气温的可能性很小。因此, 一般温度在 10 摄氏度以上对人类活动影响不大。温度字典定义如式(5)所示。

$$T = \frac{1}{1 + e^{\left(\frac{high+low}{2} + 10\right)}} \quad (5)$$

其中: *high* 是当天最高的温度; 而 *low* 是当天最低的温度; T 是处理结果。温度字典示例如图 5 所示。其中实线代表每天原始的温度, 数值对应左坐标轴; 虚线代表处理后构成的温度字典, 数值对应右坐标轴。

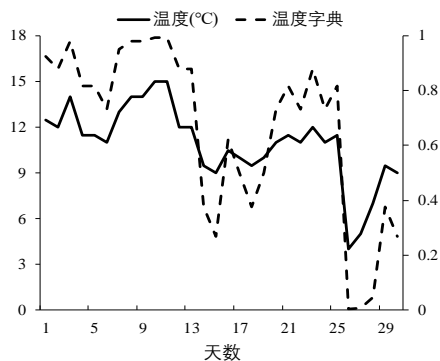


图 5 温度字典例子

Fig. 5 Example of temperature dictionary

1.3.4 天气字典

除了节假日和温度因素外, 天气也是值得考虑的重要因素。显然, 人们喜欢在天气晴朗的时候出门, 而下雨天和下雪天会减少人们的外出, 这也可能影响到商店客流量。天气情况比较复杂多样, 包括晴朗、多雨、多雪等情况。本文根据常见天气状况及其严重程度, 将天气分为以下几种不同情况, 并分别设置不同的标签, 如表 1 所示。图 6 是天气字典的例子。

表 1 天气状况分类及其标签

Table 1 Different weather conditions and their labels	
不同因素	SMAPE
W/O	0.2362
W	0.2368
T	0.2362
H	0.2088
H+T+W	0.1876

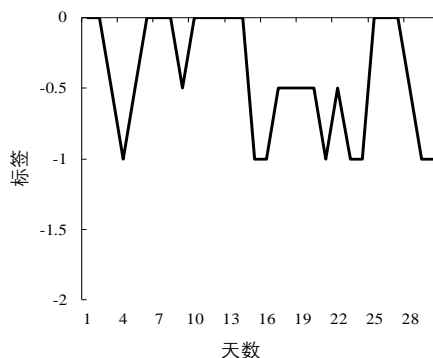


图 6 天气字典例子

Fig. 6 Example of weather dictionary

1.4 多因素稀疏编码字典

式(6)是稀疏表示过程^[17], 其目的是搜寻字典 D_f 中基向量的稀疏线性组合, 求解稀疏系数向量 α , 即使用字典 D_f 中的基向量近似地表示 y 。在本文中, 为了更准确地预测商店的客流量, 本文综合外部因素构建了多因素词典 D_f , 如式(7)所示。

$$y'' = D_f \alpha \quad (6)$$

$$D_f = [D, H, T, W] \quad (7)$$

$$\alpha = [\alpha_D, \alpha_H, \alpha_T, \alpha_W] \quad (8)$$

其中: y'' 是预处理后的商家客流量历史数据; D 是由离散余弦变换和克罗内克函数组成的字典; H 是节假日字典; T 是温度字典; W 是天气字典; α_D 、 α_H 、 α_T 和 α_W 分别是 D 字典、节假日字典、温度字典、天气字典的权重系数。

在完成过完备多因素字典 D_f 的构建之后, 本文的目的是根据式(6)求解 α , 其目标函数可以严格定义为式(9)。

$$\min \|\alpha\|_0 \quad s.t. \quad y'' = D_f \alpha \quad (9)$$

其中: α 是稀疏系数向量; $\|\alpha\|_0$ 表示 α 中非 0 的个数。上述问题的求解是一个 L_0 范数问题, 由于其非凸的、不连续的特性, 该问题的求解是一个 NP 难问题。文献[18]表明, 在满足一定的条件下, L_0 范数最小化问题可以转换成 L_1 范数最小化问题, 因此, 上述问题的求解可以转换为

$$\min \|\alpha\|_1 \quad s.t. \quad y'' = D_f \alpha \quad (10)$$

由于噪声数据的存在, 往往不能准确求解得到 α 使得式(10)完全等价。针对该问题, 通常使用一个二次惩罚函数减弱约束条件来对其进行求解, 如(11)所示。

$$\|y'' - D_f \alpha\|_2^2 < \xi \quad (11)$$

其中: ζ 是误差容忍度, 即求解得到的稀疏系数 α 使得 $\|y - D_f \alpha\|_2^2$ 的误差在给定的误差范围之内。因此, 该问题可以进一步转换成式(12) 的求解。

$$\min \frac{1}{2} \|y - D_f \alpha\|_2^2 + \lambda \|\alpha\|_1 \quad (12)$$

上述优化问题与 LASSO (Least absolute shrinkage and selection operator) [19] 的拉格朗日形式一致, 其基本思想是在回归系数绝对值之和小于某个常数的约束条件下, 最小化残差平方和, 从而可以产生一些严格等于零的回归系数, 由此来保证 α 的稀疏。本文使用 Koh 等人 [20] 提出的方法来求解 α 。求得 α 之后, α 便可以与字典 D_f 结合预测某段时间的客流量。

2 实验结果及其分析

2.1 实验设置

2.1.1 数据集

本文使用的数据集来自阿里巴巴举办的天池大数据竞赛 (口碑商家客流量预测 [21])。数据集包括全国不同城市的 2 000 家商店的历史客流量, 且其中北京、上海、杭州等移动支付较发达地区较多, 城市天气和温度数据包括每个城市的每天的情况。数据集时间范围为 2015 年 7 月 1 日至 2016 年 10 月 31 日 (除去 2015 年 12 月 12 日的数据)。

2.1.2 度量标准

本文使用 2 000 个商店的平均误差作为每组实验的误差, 以此来衡量预测方法的好坏。由于某些商家的真实客流量比较小, 甚至为 0, 为了减少这类商家对总体误差带来的影响。在本文中, 本文使用对称平均绝对百分比误差 (symmetric mean absolute percentage error, SMAPE) 作为每组实验的预测误差的计算方法, 如式(13) 所示。

$$SMAPE = \frac{2}{nT} \sum_{i=1}^n \sum_{t=1}^T \frac{|c_{it} - c_{it}^e|}{|c_{it} + c_{it}^e|} \quad (13)$$

其中: n 是商店的数量; T 是预测的天数; c_{it} 是第 i 天的预测值; c_{it}^e 是第 i 天的实际值; $SMAPE$ 是最终误差。

2.2 不同外部因素的实验结果

在该实验中, 本文对没有因素、增加一个因素、增加所有因素的不同情况进行实验, 并分别计算平均预测误差。实验结果如表 2 所示。其中: W/O 是不加因素; W、T、H 分别是天气因素、温度因素和节假日因素。显然, 与不加因素相比, 加上节假日因素, 温度因素和天气因素后, 预测精度提高了 4.86%。在只增加节假日因素的情况下, 预测的准确性在一定程度上有所提高, 为 2.74%。然而只在增加天气或温度因素时, 预测精度几乎不变。从这个角度来看, 节假日因素比天气因素和温度因素的影响更大。

表 2 不同因素的预测误差

不同因素	SMAPE
W/O	0.2362
W	0.2368
T	0.2362
H	0.2088
H+T+W	0.1876

图 7 是 2 000 个商家三种因素权重的分布图, 分别对应式(8)中的 α_H 、 α_T 、 α_W 。总体上, 节假日因素的比例最大, 温度其次, 天气最小, 意味着节假日因素的影响最大, 温度和天气因素影响相对较小, 与上述实验结果对应。同时, 也可以看到某些商家的天气因素权重很大, 说明天气对某些商家

的影响特别大。

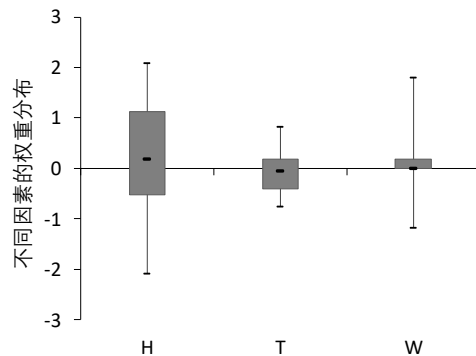


图 7 不同因素的权重分布

Fig. 7 Distribution of different factors weight

综上所述, 从总体上而言, 仅仅使用节假日因素对预测精度有明显的提高, 单独使用温度或者天气因素对于预测精度的影响很小。实验结果表明, 外部因素对于商家客流量确实有较大影响。

2.3 不同商家类型的实验结果

由于店铺种类不同, 预测模型在这些商店的预测效果可能会有所不同, 不同店铺对不同因素敏感度也不同。在 2 000 家的商店里, 有 579 间超市、639 间餐厅和 782 间小吃店。本文在这些不同类型的商店上进行实验, 并分析外部因素对不同商店的影响。

如图 8 所示, 其中 W/O 表示不附加任何外部因素, W、T、H 分别表示附加天气、温度和节假日因素。预测模型在超市的预测结果最好, 小吃店最差。与超市和餐厅相比, 小吃店的规模要小得多, 因此, 日常客流量相对不稳定, 导致小吃店的预测效果最差。在增加所有外部因素后, 超市预测精度提高了 3.94%, 餐厅预测精度提高 4.88%, 而小吃店的预测精度提高了 5.46%。可能的原因是因为在平时生活中人们不得不去超市购买生活必需品, 而餐厅和小吃店的消费相对来说没有那么重要, 导致外部因素对超市的影响相对较小。因此, 小吃店相比超市和餐厅对外部因素更为敏感。节假日仍然在三种因素中影响最大, 这与 2.2 节 所得的结论一致。

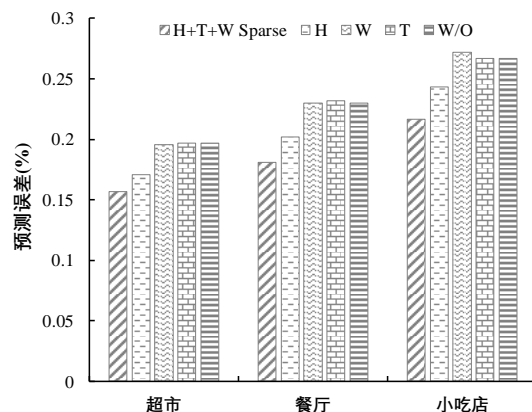


图 8 不同商店类型的预测结果

Fig. 8 Prediction error of different shops

图 9 是商家数量在不同因素权重范围的分布情况。其中图 9(a)(b)和(c)分别是在节假日、温度、天气因素下, 商家数量在不同权重范围的百分比。如图 9(a)所示, 总体上, 由于三种类型商家数量的分布偏向于权重为正, 所以节假日因素对于客流量的影响偏向于积极。同理, 对于温度因素, 如图

9(b)所示, 温度因素对于三种类型商家客流量的影响偏向消极。同时可以看到, 超市在节假日因素和温度因素权重值为 0 附近的数量的百分比显著大于餐厅和小吃店, 当权重值变大或变小时, 超市数量的百分比普遍小于这两种类型商家, 这意味着相比餐厅和小吃店, 节假日因素和温度因素对超市

客流量产生的影响要小。商家数量在天气因素不同权重的分布如图 9(c)所示, 显然三种不同类型商家的分布几乎没有区别, 特别是当权重值在 0 附近, 三种商家数量的百分比都达到了 70%以上, 并随着权重值的改变, 商家数量迅速下降, 这表明天气因素对于三种不同类型的商家客流量的影响很小。

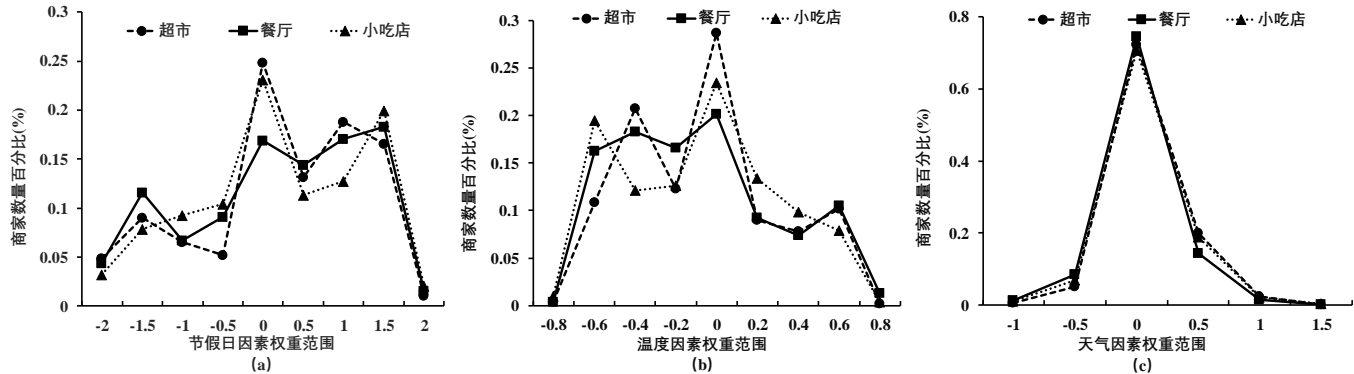


图 9 商家数量在不同因素权重范围的分布情况
Fig. 9 Distribution of shops in different factors weight

2.4 不同预测模型的实验结果

在该实验中, 本文选取其他四种预测方法来和多因素稀疏回归模型 (即 H+T+W Sparse 模型) 进行对比, 其中所有预测方法的训练数据长度为 30 d, 测试数据长度为随后 14 d。第一种方法为 ARIMA 模型, 本文采用参数估计^[3]的方式, 对每个商家数据经过差分、定阶来确定 ARIMA 模型的参数; 第二种是采用取均值的方法, 即预测的客流量全为训练数据客流量的均值; 第三种是不附加任何外部因素的 Sparse Regression 模型; 第四种是循环神经网络 (recurrent neural network, RNN) ^[22], 其中网络的隐藏层层数为 20, 学习率为 0.006。实验结果如表 3 所示。

表 3 不同模型的预测误差

Table 3 Prediction error of different models

预测模型	SMAPE
RNN	0.2628
MEAN	0.2490
Sparse Regression	0.2362
ARIMA	0.2154
H+T+W Sparse	0.1876

根据上述实验结果, 多因素稀疏回归模型的效果要明显好于其他四种预测方法, 其中 RNN 模型预测效果最差。不同模型预测误差的 CDF 图如图 10 所示。从 CDF 图分析, 五种方法在误差为 0.1 和 0.6 以下时, 商家数量区别不明显; 但是当误差在 0.1~0.5 间, 特别是当误差为 0.2 和 0.3 时, H+T+W Sparse 相比其他方法, 在商家数量上有着 10%~50% 的提升。显然, 在 Sparse Regression 模型的基础上增加节假日、温度和天气因素后, 预测效果得到显著的提高。

预测模型的时间复杂度也是衡量模型性能的重要指标之一。在本文中, 由于有 2 000 个商家, 所以本文计算所有商家在不同模型上的训练时间以及预测时间的均值来对比上述模型的时间复杂度。取均值的预测方法不能算真正意义上的预测模型, 因此不考虑在内。实验结果如表 4 所示。

其中 RNN 实验平台为 GPU(GTX1080Ti), 其余三个模型的实验平台是 CPU(i7-4790), 内存均为 8 GB。从表 4 中可以看出, 稀疏回归模型和多因素稀疏回归模型的训练时间以及预测时间都要远远少于其他两个模型。而多因素稀疏回归模型与稀疏回归模型的时间复杂度几乎没有差距, 考虑到预测的准确度, 显然, 多因素时间预测模型的性能更加优越。

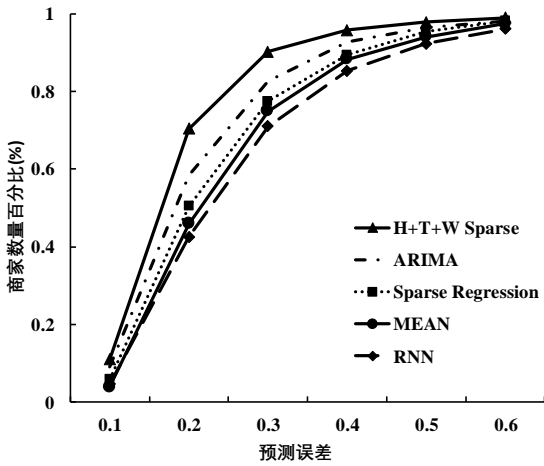


图 10 不同模型预测误差的 CDF 图
Fig. 10 CDF of different model prediction error

表 4 不同模型的时间复杂度

Table 4 Time complexity of different models

预测模型	训练时间/ms	预测时间/ms
RNN	506.1	433.3
Sparse Regression	4.7	0.03
ARIMA	524.7	50.2
H+T+W Sparse	5.1	0.03

2.5 不同训练长度的实验结果

训练长度对建立预测模型至关重要, 不同的训练长度导致预测结果差异很大。一方面, 过长的训练长度包含过多的历史数据, 会对预测结果造成负面影响; 另一方面, 训练长度太短会导致无法提取足够的特征, 不能准确表示客流量的趋势。在本文中, 商店的客流量随着时间的推移而变化, 对于同一个商家而言, 相隔几个月的客流量可能都会有很大的不同。因此, 本文以多因素稀疏回归方法为预测模型, 选择了 20 ~100 的训练长度, 并计算每种训练长度的误差。实验结果如图 11 所示。当训练长度为 30 时, 预测结果最好。

3 结束语

本文提出了一个结合节假日、温度和天气因素的多因素稀疏回归预测模型。研究过程中, 本文分析了外部因素可能带来的影响, 并验证这些因素对客流量预测精度的影响。同

时, 通过在不同商家类型上进行预测实验, 分析了不同商家对外界因素的敏感性。实验结果表明, 综合多个外部因素建立的预测模型显著好于无因素模型, 同时发现部分外界因素与商家客流量有着密切的关系。本文研究成果对于提高预测客流量的准确度具有一定的现实意义。

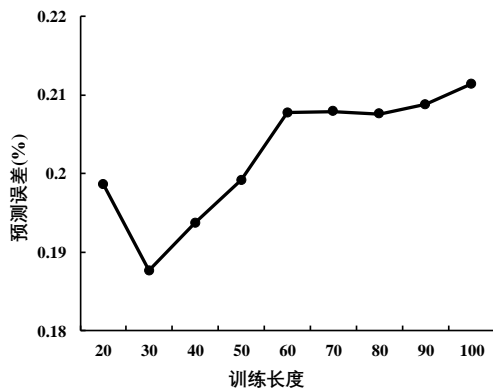


图 11 不同训练长度的预测误差

Fig. 11 Prediction error of different training length

参考文献:

- [1] Bippig N. Best practice in inventory management - 2nd edn [J]. International Journal of Production Research, 2003, 41 (15): 3645-3645.
- [2] 龚其国, 黄文辉. 供应链管理中集中库存研究综述与展望 [J]. 管理评论, 2017, 29 (11): 206-215. (Gong Qiguo, Huang Wenhui. review and prospects of inventory pooling in supply chain management [J]. Business Review, 2017, 29 (11): 206-215.)
- [3] Box G E P, Jenkins G M. Time series analysis: forecasting and control [J]. Journal of Time, 2010, 31 (4): 303-303.
- [4] Ramos P, Santos N, Rui R. Performance of state space and ARIMA models for consumer retail sales forecasting [J]. Robotics and Computer Integrated Manufacturing, 2015, 34: 151-163.
- [5] 刘建军, 廖闻剑, 彭艳兵. 两种时间序列模型在客流量预测上的比较 [J]. 计算机工程与应用, 2016, 52 (9): 228-232. (Liu Jianjun, Liao Wenjian, Peng Yanbing. Comparison between two kinds of time series models for forecasting passenger flow [J]. Computer Engineering and Applications, 2016, 52 (9): 228-232.)
- [6] Liu L M, Bhattacharyya S, Sclove S L, et al. Data mining on time series: an illustration using fast-food restaurant franchise data [J]. Computational Statistics & Data Analysis, 2001, 37 (4): 455-476.
- [7] 杨本臣, 王翠琴. P2P 中基于云模型和灰色系统理论的信任机制研究 [J]. 计算机应用研究, 2016, 33 (1): 276-280. (Yang Benchen, Wang Cuiqin. Research on trust mechanism based on cloud model and gray system theory for P2P network [J]. Application Research of Computers. 2016, 33 (1): 276-280.)
- [8] 刘卫校. 基于离散灰色预测模型与神经网络混合智能模型的时尚销售预测 [J]. 计算机应用, 2016, 36 (12): 3378-3384. (Liu Weixiao. Hybrid intelligent model for fashion sales forecasting based on discrete grey forecasting model and artificial neural network [J]. Journal of Computer Applications, 2016, 36 (12): 3378-3384.)
- [9] İşlek İ, Ögüdücü Ş G. A retail demand forecasting model based on data mining techniques [C]//Proc of the 24th IEEE International Symposium on Industrial Electronics. Piscataway, NJ: IEEE Press, 2015: 55-60.
- [10] Schneider M J, Gupta S. Forecasting sales of new and existing products using consumer reviews: a random projections approach [J]. International Journal of Forecasting, 2016, 32 (2): 243-256.
- [11] Merino M, Ramirez-Nafarrate A. Estimation of retail sales under competitive location in Mexico [J]. Journal of Business Research, 2015, 69 (2): 445-451.
- [12] Luo J. Integrating the huff model and floating catchment area methods to analyze spatial access to healthcare services [J]. Trans in Gis, 2014, 18 (3): 436-448.
- [13] Dufo-López R, Pérez-Cebollada E, Bernal-Agustín J L, et al. Optimisation of energy supply at off-grid healthcare facilities using Monte Carlo simulation [J]. Energy Conversion & Management, 2016, 113: 321-330.
- [14] Cao L, Jin L, Tao H, et al. Multi-focus image fusion based on spatial frequency in discrete cosine transform domain [J]. IEEE Signal Processing Letters, 2015, 22 (2): 220-224.
- [15] Qayyum A, Malik A S, Naufal M, et al. Designing of overcomplete dictionaries based on DCT and DWT [C]// Proc of IEEE Student Symposium in Biomedical Engineering & Sciences. Piscataway, NJ: IEEE Press, 2015: 134-139.
- [16] Shi Y, Gao Y, Yang Y, et al. Multimodal sparse representation-based classification for lung needle biopsy images [J]. IEEE Trans on Biomedical Engineering, 2013, 60 (10): 2675-2685.
- [17] 邓欣, 高峰星, 米建勋, 等. 基于稀疏表示的脑电 (EEG) 情感分类 [J/OL]. 计算机应用研究, 2019 (4): 1-2 [2018-12-16]. <http://kns.cnki.net/kcms/detail/51.1196.TP.20180209.1230.174.html>. (Deng Xin, Gao Fengxing, Mi Jianxun, et al. Classifying emotional EEG using sparse representation method [J/OL]. Application Research of Computers, 2019 (4): 1-2 [2018-12-16]. <http://kns.cnki.net/kcms/detail/51.1196.TP.20180209.1230.174.html>.)
- [18] Donoho D L. Compressed sensing [J]. IEEE Trans on Information Theory, 2006, 52 (4): 1289-1306.
- [19] Tibshirani R. Regression shrinkage and selection via the lasso: a retrospective [J]. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 2011, 73 (3): 273-282.
- [20] Kim S J, Koh K, Lustig M, et al. An interior-point method for large-scale l1-regularized least squares [J]. IEEE Journal of Selected Topics in Signal Processing, 2007, 1 (4): 606-617.
- [21] IJCAI-17 口碑商家客流量预测 [EB/OL]. [2018-08-14]. <https://tianchi.aliyun.com/competition/introduction.htm?spm=5176.11409391.333.4.58ee49feRuZHFB&raceId=231591>.
- [22] Balluff S, Bendfeld J, Krauter S. Meteorological data forecast using RNN [J]. International Journal of Grid and High Performance Computing, 2017, 9 (1): 61-74.